

## Seminario

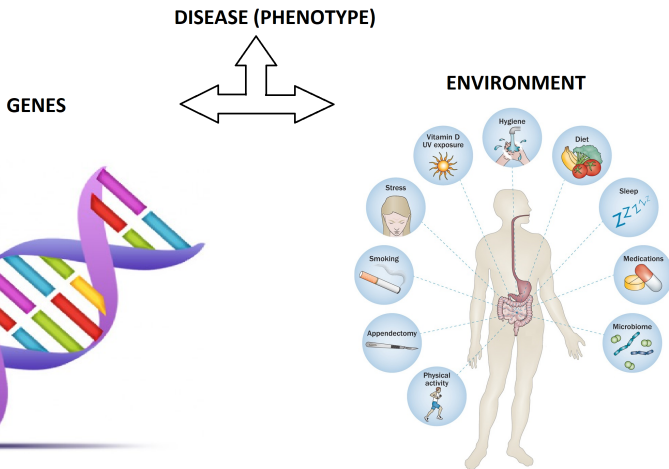
Testing gene-environment interaction in generalized linear mixed models with family data

20 de noviembre de 2017

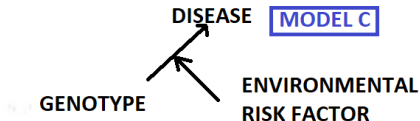
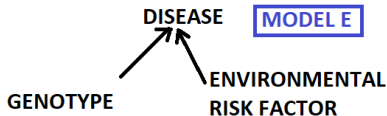
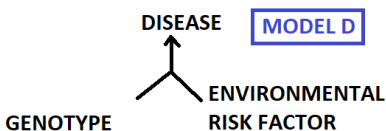
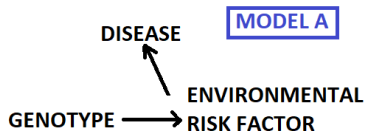
Sección 1

# Introduction

Gene, Environment and Disease

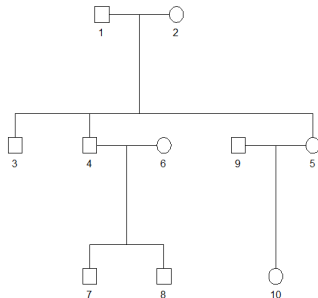


Gene-environment interaction



## Family data and kinship matrix

**PEDIGREE**



**KINSHIP MATRIX**

Subject	1	2	3	4	5	6	7	8	9	10
1	1.00	0.00	0.50	0.50	0.50	0.0	0.250	0.250	0.0	0.250
2	0.00	1.00	0.50	0.50	0.50	0.0	0.250	0.250	0.0	0.250
3	0.50	0.50	1.00	0.50	0.50	0.0	0.250	0.250	0.0	0.250
4	0.50	0.50	0.50	1.00	0.50	0.0	0.500	0.500	0.0	0.250
5	0.50	0.50	0.50	0.50	1.00	0.0	0.250	0.250	0.0	0.500
6	0.00	0.00	0.00	0.00	0.00	1.0	0.500	0.500	0.0	0.000
7	0.25	0.25	0.25	0.50	0.25	0.5	1.000	0.500	0.0	0.125
8	0.25	0.25	0.25	0.50	0.25	0.5	0.500	1.000	0.0	0.125
9	0.00	0.00	0.00	0.00	0.00	0.0	0.000	0.000	1.0	0.500
10	0.25	0.25	0.25	0.25	0.50	0.0	0.125	0.125	0.5	1.000

## Notation

- Consider a sample of  $N$  independent families.
- $n_i$ : number of members in the  $i$ th family ( $i = 1, \dots, N$ ).
- $Y_{ij}$ : reponse variable for the phenotype (discrete or continuous).
- $\mathbf{X}_{ij} = (X_{ij}^1, \dots, X_{ij}^p)^T$ :  $p$  non-environmental covariates.
- $\mathbf{G}_{ij} = (G_{ij}^1, \dots, G_{ij}^q)^T$ :  $q$  observed genotypes at certain targeted genetic marker loci.
- $E_{ij}$ : some environmental exposure factor of interest.
- $\mathbf{S}_{ij} = (E_{ij}G_{ij}^1, \dots, E_{ij}G_{ij}^q)^T$ : GE interaction.

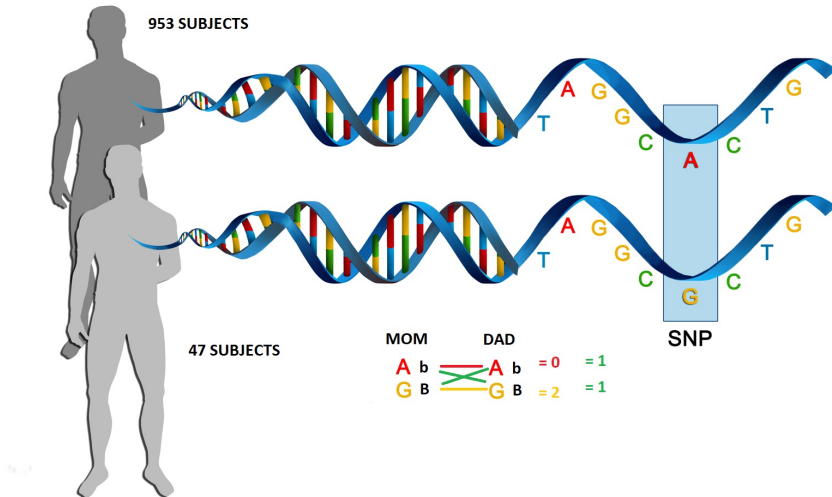
Observed genotypes -  $G_{ij}$  -

Each column in  $G_{ij}$  is a Single Nucleotide Polymorphism (SNP). The genetic information is represented according to the following codification:

$$G_{ij}^k = \begin{cases} 2, & \text{subject } j \text{ at } i\text{th family is homozygous BB} \\ 1, & \text{subject } j \text{ at } i\text{th family is heterozygous Bb or bB} \\ 0, & \text{subject } j \text{ at } i\text{th family is homozygous bb,} \end{cases}$$

with  $k = 1, \dots, q$ . B and b represent the dominant and recessive alleles, respectively. In addition, B is the allele that occurs at minor frequency.

Observed genotypes -  $G_{ij}$  -





Sección 2

# Model

Generalized linear mixed model (GLMM) for Subject  $j$  at  $i$ th family

$$g[E(Y_{ij}|\alpha_{ij})] = \mathbf{X}_{ij}^T \beta_1 + E_{ij} \beta_2 + \mathbf{G}_{ij}^T \boldsymbol{\theta} + \mathbf{S}_{ij}^T \boldsymbol{\gamma} + \alpha_{ij}, \quad (1)$$

$$\text{Var}(Y_{ij}|\alpha_{ij}) = \phi \omega_{ij}^{-1} \nu[E(Y_{ij}|\alpha_{ij})],$$

$$\boldsymbol{\alpha}_i = (\alpha_{i1}, \dots, \alpha_{in_i})^T \sim N(\mathbf{0}, 2\sigma^2 \boldsymbol{\Phi}_i),$$

where,

- $g(\cdot)$ : monotone known function.
- $Y_{ij}|\alpha_{ij}$  follows a distribution in the exponential family.
- $\nu(\cdot)$ : known function.
- $\phi$ : a scale parameter that may be known or may need to be estimated.
- $\omega_{ij}$ : known weights (commonly equal to 1).
- $\boldsymbol{\Phi}_i$ : the kinship matrix and  $\sigma^2$  is a parameter to be estimated.

## Examples of link functions

Family	Link	$g(\cdot)$	Trait	$\nu(\cdot)$
Binomial	Logit	$\ln\left(\frac{\mu}{1-\mu}\right)$	Binary	$\mu(1-\mu)$
Gaussian	Identity	$\mu$	Continuous	$\phi$
Gamma	Inverse	$1/\mu$	Continuous	$\phi\mu^2$
Inverse.gaussian	Inverse squared	$1/\mu^2$	Continuous	$\phi\mu^3$
Poisson	Log	$\ln(\mu)$	Count	$\mu$
Quasi	Identity	$\mu$	Continuous	$\phi$
Quasibinomial	Logit	$\ln\left(\frac{\mu}{1-\mu}\right)$	Binary	$\phi\mu(1-\mu)$
Quasipoisson	Log	$\ln(\mu)$	Count	$\phi\mu$

GLMM for  $i$ th Family

$$g(\boldsymbol{\mu}_i^b) = \mathbf{X}_i\boldsymbol{\beta}_1 + \mathbf{E}_i\boldsymbol{\beta}_2 + \mathbf{G}_i\boldsymbol{\theta} + \mathbf{S}_i\boldsymbol{\gamma} + \mathbf{K}_i\mathbf{b}_i, \quad (2)$$

where,

- $2\Phi_i = \mathbf{K}_i\mathbf{K}_i^T$  (Cholesky decomposition).
- $\boldsymbol{\alpha}_i = \mathbf{K}_i\mathbf{b}_i$ , with  $\mathbf{b}_i \sim N(\mathbf{0}, \sigma^2\mathbf{I}_{n_i})$  and  $\mathbf{I}_{n_i}$ : identity matrix.
- $\boldsymbol{\mu}_i^b = E(\mathbf{Y}_i|\mathbf{b}_i)$ .
- $g(\boldsymbol{\mu}_i^b) = (g(\mu_{i1}^b), \dots, g(\mu_{in_i}^b))^T$ .
- $\mathbf{X}_i = [\mathbf{X}_{i1} \dots \mathbf{X}_{in_i}]^T$ .
- $\mathbf{G}_i = [\mathbf{G}_{i1} \dots \mathbf{G}_{in_i}]^T$ .
- $\mathbf{E}_i = (E_{i1}, \dots, E_{in_i})^T$ .
- $\mathbf{S}_i = [\mathbf{S}_{i1} \dots \mathbf{S}_{in_i}]^T$ .

General GLMM

$$g(\boldsymbol{\mu}^b) = \tilde{\mathbf{X}}\boldsymbol{\beta} + \mathbf{G}\boldsymbol{\theta} + \mathbf{S}\boldsymbol{\gamma} + \mathbf{K}\mathbf{b}, \quad (3)$$

where

- $\boldsymbol{\mu}^b = E(\mathbf{Y}|\mathbf{b})$ .
- $\mathbf{X} = [\mathbf{X}_1 \dots \mathbf{X}_N]^T$ .
- $\mathbf{G} = [\mathbf{G}_1 \dots \mathbf{G}_N]^T$ .
- $\mathbf{E} = (E_1, \dots, E_N)^T$ .
- $\mathbf{S} = [\mathbf{S}_1 \dots \mathbf{S}_N]^T$ .
- $\mathbf{K} = \text{diag}\{\mathbf{K}_1 \dots \mathbf{K}_N\}$ .
- $\mathbf{b} = [\mathbf{b}_1 \dots \mathbf{b}_N]^T$ .
- $\tilde{\mathbf{X}} = [\mathbf{X} \ \mathbf{E}]^T$ .
- $\boldsymbol{\beta} = (\boldsymbol{\beta}_1^T, \boldsymbol{\beta}_2^T)^T$ .

We are interested in testing the hypothesis  $H_0 : \boldsymbol{\gamma} = \mathbf{0}$ .

### Important facts

- If  $\gamma$  is treated as a fixed vector and the null hypothesis is tested with a  $q$  degrees of freedom score test, it can result in loss of power (Lin *et. al.*, 2013).



### Important facts

- If  $\gamma$  is treated as a fixed vector and the null hypothesis is tested with a  $q$  degrees of freedom score test, it can result in loss of power (Lin *et. al.*, 2013).
- Another common strategy is to use a single SNP at time to test GE interaction.

### Important facts

- If  $\gamma$  is treated as a fixed vector and the null hypothesis is tested with a  $q$  degrees of freedom score test, it can result in loss of power (Lin *et. al.*, 2013).
- Another common strategy is to use a single SNP at time to test GE interaction.
- Assume  $\gamma$  as a random vector following a multivariate normal distribution  $N(\mathbf{0}, \tau \mathbf{I}_q)$  and to test the equivalent null hypothesis  $H_0 : \tau = 0$ .



### Important facts

- If  $\gamma$  is treated as a fixed vector and the null hypothesis is tested with a  $q$  degrees of freedom score test, it can result in loss of power (Lin *et. al.*, 2013).
- Another common strategy is to use a single SNP at time to test GE interaction.
- Assume  $\gamma$  as a random vector following a multivariate normal distribution  $N(\mathbf{0}, \tau \mathbf{I}_q)$  and to test the equivalent null hypothesis  $H_0 : \tau = 0$ .

$$g(\mu^{b,\gamma}) = \tilde{\mathbf{X}}\beta + \mathbf{G}\theta + \underbrace{\mathbf{S}\gamma + \mathbf{K}b}_{\text{random}}$$

### Important facts

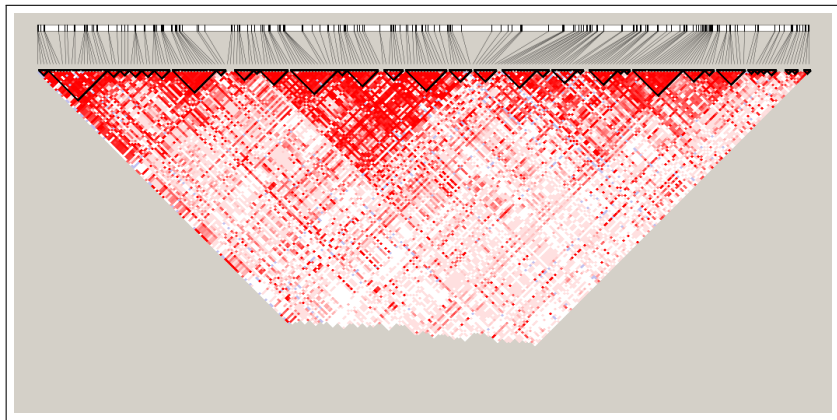
- If  $\gamma$  is treated as a fixed vector and the null hypothesis is tested with a  $q$  degrees of freedom score test, it can result in loss of power (Lin *et. al.*, 2013).
- Another common strategy is to use a single SNP at time to test GE interaction.
- Assume  $\gamma$  as a random vector following a multivariate normal distribution  $N(\mathbf{0}, \tau \mathbf{I}_q)$  and to test the equivalent null hypothesis  $H_0 : \tau = 0$ .

$$g(\mu^b) = \tilde{\mathbf{X}}\beta + \mathbf{G}\theta + \underbrace{\mathbf{S}\gamma}_{\text{random}} + \underbrace{\mathbf{K}b}_{\text{random}}$$

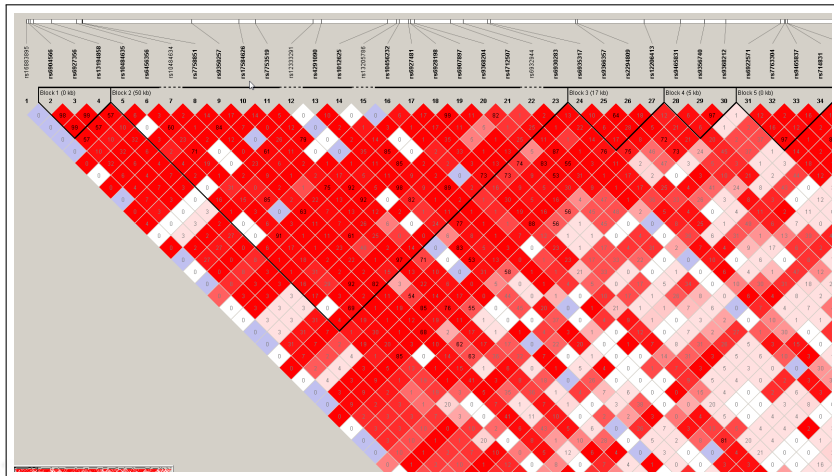
## Linkage disequilibrium (LD) -PPARG-



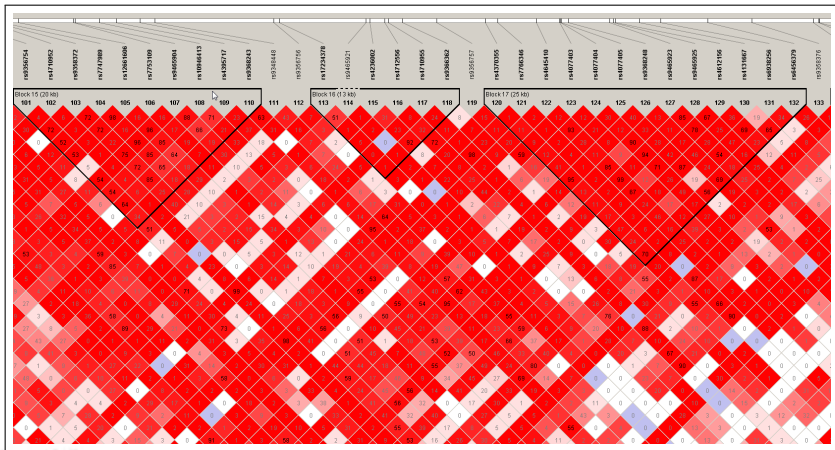
## Linkage disequilibrium (LD) -CDKAL1-



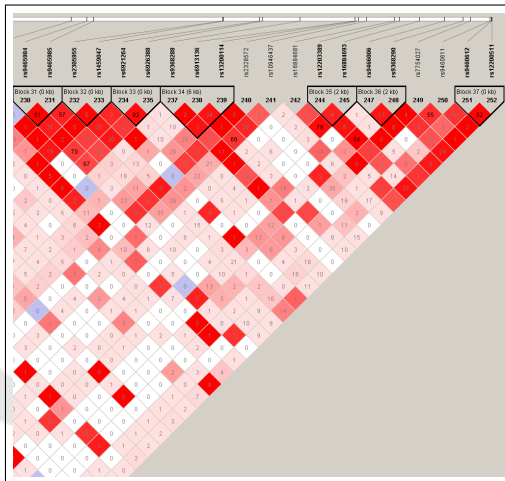
# Linkage disequilibrium (LD) -CDKAL1-



## Linkage disequilibrium (LD) -CDKAL1-



## Linkage disequilibrium (LD) -CDKAL1-



## Remedial considerations to face LD

- Given the high correlation among the SNPs in  $\mathbf{G}$ , Lin *et. al.* (2013) proposed (for independent subjects) to penalize the estimation of parameter  $\theta$  by using ridge regression and introducing a penalization term  $\lambda$  in the quasi-likelihood function. The tuning parameter  $\lambda$  is selected using generalized cross validation (Fu, 2005).



## Remedial considerations to face LD

- Given the high correlation among the SNPs in  $\mathbf{G}$ , Lin *et. al.* (2013) proposed (for independent subjects) to penalize the estimation of parameter  $\theta$  by using ridge regression and introducing a penalization term  $\lambda$  in the quasi-likelihood function. The tuning parameter  $\lambda$  is selected using generalized cross validation (Fu, 2005).
- Shen *et. al.* (2013) proposed for generalized linear models (GLM) that ridge regression is equivalent to assume the penalized parameters as independent random variables.

## Remedial considerations to face LD

It is also equivalent for GLMM and assuming  $\boldsymbol{\theta} \sim N(\mathbf{0}, \sigma_{\theta}^2 \mathbf{I}_q)$ , it is possible to show that  $\lambda = \phi / \sigma_{\theta}^2$ .

## Remedial considerations to face LD

It is also equivalent for GLMM and assuming  $\boldsymbol{\theta} \sim N(\mathbf{0}, \sigma_{\theta}^2 \mathbf{I}_q)$ , it is possible to show that  $\lambda = \phi / \sigma_{\theta}^2$ .

$$g\left(\boldsymbol{\mu}^{b,\gamma,\theta}\right) = \tilde{\mathbf{X}}\boldsymbol{\beta} + \underbrace{\mathbf{G}\boldsymbol{\theta} + \mathbf{S}\boldsymbol{\gamma} + \mathbf{K}\mathbf{b}}_{\text{random}}$$

Remedial considerations to face LD

It is also equivalent for GLMM and assuming  $\theta \sim N(\mathbf{0}, \sigma_\theta^2 \mathbf{I}_q)$ , it is possible to show that  $\lambda = \phi / \sigma_\theta^2$ .

$$g(\mu^{b,\theta}) = \tilde{\mathbf{X}}\beta + \underbrace{\mathbf{G}\theta}_{\text{random}} + \mathbf{S}\gamma + \underbrace{\mathbf{K}b}_{\text{random}}$$

Sección 3

# Null model estimation

Null model estimation

$$g\left(\mu^{d_1, d_2}\right) = \tilde{\mathbf{X}}\beta + \mathbf{d}_1 + \mathbf{d}_2$$

with  $\mathbf{d}_1 = \mathbf{G}\theta$ ,  $\mathbf{d}_2 = \mathbf{K}b$ .



Null model estimation

$$g(\boldsymbol{\mu}^{d_1, d_2}) = \tilde{\mathbf{X}}\boldsymbol{\beta} + \mathbf{d}_1 + \mathbf{d}_2$$

with  $\mathbf{d}_1 = \mathbf{G}\boldsymbol{\theta}$ ,  $\mathbf{d}_2 = \mathbf{K}\mathbf{b}$ . Breslow and Clayton (1993) proposed a Fisher scoring solution that may be expressed as the iterative solution to the system

$$\begin{bmatrix} \tilde{\mathbf{X}}^T \mathbf{W} \tilde{\mathbf{X}} & \tilde{\mathbf{X}}^T \mathbf{W} & \tilde{\mathbf{X}}^T \mathbf{W} \\ \mathbf{W} \tilde{\mathbf{X}} & \frac{\phi}{\sigma^2} (\mathbf{G}\mathbf{G}^T)^{-1} + \mathbf{W} & \mathbf{W} \\ \mathbf{W} \tilde{\mathbf{X}} & \mathbf{W} & \frac{\phi}{\sigma^2} (\mathbf{K}\mathbf{K}^T)^{-1} + \mathbf{W} \end{bmatrix} \begin{pmatrix} \boldsymbol{\beta} \\ \mathbf{d}_1 \\ \mathbf{d}_2 \end{pmatrix} = \begin{pmatrix} \tilde{\mathbf{X}}^T \mathbf{W} \tilde{\mathbf{Y}} \\ \mathbf{W} \tilde{\mathbf{Y}} \\ \mathbf{W} \tilde{\mathbf{Y}} \end{pmatrix}$$

- $\tilde{\mathbf{Y}} = \tilde{\mathbf{X}}\boldsymbol{\beta} + \mathbf{d}_1 + \mathbf{d}_2 + \boldsymbol{\varepsilon}$ : working vector;
- $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \phi \mathbf{W}^{-1})$ ;
- $\mathbf{W} = \text{diag} \left\{ \omega_{ij} / \left[ \nu(\mu_{ij}^d) g'(\mu_{ij}^d)^2 \right] \right\}$ .

Null model estimation

$$\begin{aligned}\hat{\beta} &= (\tilde{\mathbf{X}}^T \Sigma^{-1} \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \Sigma^{-1} \tilde{\mathbf{Y}}, \\ \begin{pmatrix} \hat{d}_1 \\ \hat{d}_2 \end{pmatrix} &= \begin{pmatrix} \sigma_{\theta}^2 (\mathbf{G}\mathbf{G}^T) \Sigma^{-1} (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\hat{\beta}) \\ \sigma^2 (\mathbf{K}\mathbf{K}^T) \Sigma^{-1} (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\hat{\beta}) \end{pmatrix},\end{aligned}$$

with  $\Sigma = \sigma_{\theta}^2 \mathbf{G}\mathbf{G}^T + \sigma^2 \mathbf{K}\mathbf{K}^T + \phi \mathbf{W}^{-1}$ .



Sección 4

# GE interaction test

## Testing GE interaction

Following the approach developed by Zhang and Lin (2003), we propose as score statistic the quadratic form

$$U_{\tau} = U_{\tau}(\hat{\beta}, \hat{\pi}) = \frac{1}{2} \left\{ (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\beta)^T \Sigma^{-1} \mathbf{S} \mathbf{S}^T \Sigma^{-1} (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\beta) \right\} \Big|_{\hat{\beta}, \hat{\pi}}.$$

with  $\hat{\pi}$  is the estimator of  $\pi = (\sigma_{\theta}^2, \sigma^2, \phi)^T$ .

## Testing GE interaction

Following the approach developed by Zhang and Lin (2003), we propose as score statistic the quadratic form

$$U_{\tau} = U_{\tau}(\hat{\beta}, \hat{\pi}) = \frac{1}{2} \left\{ \left( \tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\beta \right)^T \Sigma^{-1} \mathbf{S} \mathbf{S}^T \Sigma^{-1} \left( \tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\beta \right) \right\} \Big|_{\hat{\beta}, \hat{\pi}}.$$

with  $\hat{\pi}$  is the estimator of  $\pi = (\sigma_{\theta}^2, \sigma^2, \phi)^T$ . To correct for bias, we use the restricted maximum likelihood (REML) estimators (Breslow and Clayton, 1993) in the GLMM framework to obtain  $\hat{\beta}$  and  $\hat{\pi}$  under the null hypothesis.

## Testing GE interaction

Zhang and Lin (2003) showed that under  $H_0 : \tau = 0$ ,  $U_\tau$  follows approximately a mixture of one degree of freedom independent chi-square distributions. However, for computational ease, we use the Satterthwaite method (Satterthwaite, 1941) to approximate the distribution of  $U_\tau$  by a scaled chi-square distribution  $\kappa\chi_\xi^2$ .

## Testing GE interaction

Zhang and Lin (2003) showed that under  $H_0 : \tau = 0$ ,  $U_\tau$  follows approximately a mixture of one degree of freedom independent chi-square distributions. However, for computational ease, we use the Satterthwaite method (Satterthwaite, 1941) to approximate the distribution of  $U_\tau$  by a scaled chi-square distribution  $\kappa\chi_\xi^2$ .

When REML estimates are used to calculate  $U_\tau$  showed that the mean and variance of  $U_\tau$  can be approximated, respectively, by

$$\text{tr}(\mathbf{PSS}^T) \Big|_{\hat{\pi}} \quad \text{and} \quad \mathcal{I}_\tau = \frac{1}{2} \left\{ \text{tr}(\mathbf{PSS}^T \mathbf{PSS}^T) - \mathbf{J}^T \mathbf{M}^{-1} \mathbf{J} \right\} \Big|_{\hat{\beta}, \hat{\pi}}$$

Testing GE interaction

$$J = \left( \text{tr}[PSS^T PGG^T] \quad \text{tr}[PSS^T PKK^T] \mid \text{tr}[PSS^T PW^{-1}] \right)^T$$

$$M = \begin{bmatrix} \text{tr}[PGG^T PGG^T] & \text{tr}[PGG^T PKK^T] & \text{tr}[PGG^T PW^{-1}] \\ \text{tr}[PKK^T PG^T G] & \text{tr}[PKK^T PKK^T] & \text{tr}[PKK^T PW^{-1}] \\ \text{tr}[PW^{-1} PG^T G] & \text{tr}[PW^{-1} PKK^T] & \text{tr}[PW^{-1} PW^{-1}] \end{bmatrix},$$

and  $P = \Sigma^{-1} - \Sigma^{-1} \tilde{X} \left( \tilde{X}^T \Sigma^{-1} \tilde{X} \right)^{-1} \tilde{X}^T \Sigma^{-1}.$

## Testing GE interaction

Since the mean and variance of  $\kappa\chi_{\xi}^2$  are given by  $\kappa\xi$  and  $2\kappa^2\xi$ , respectively, we obtain the equations  $\text{tr}(\hat{\mathbf{P}}\mathbf{S}\mathbf{S}^T) = \kappa\xi$  and  $\mathcal{I}_{\tau} = 2\kappa^2\xi$ , where  $\hat{\mathbf{P}}$  denotes the matrix  $\mathbf{P}$  evaluated in  $\hat{\boldsymbol{\pi}}$ . By solving these equations, we demonstrate that

$$\kappa = \mathcal{I}_{\tau} / [2 \text{tr}(\hat{\mathbf{P}}\mathbf{S}\mathbf{S}^T)]$$

and

$$\xi = 2 \left[ \text{tr}(\hat{\mathbf{P}}\mathbf{S}\mathbf{S}^T) \right]^2 / \mathcal{I}_{\tau}.$$

## Testing GE interaction

Since the mean and variance of  $\kappa\chi_{\xi}^2$  are given by  $\kappa\xi$  and  $2\kappa^2\xi$ , respectively, we obtain the equations  $\text{tr}(\hat{\mathbf{P}}\mathbf{S}\mathbf{S}^T) = \kappa\xi$  and  $\mathcal{I}_{\tau} = 2\kappa^2\xi$ , where  $\hat{\mathbf{P}}$  denotes the matrix  $\mathbf{P}$  evaluated in  $\hat{\pi}$ . By solving these equations, we demonstrate that

$$\kappa = \mathcal{I}_{\tau} / [2 \text{tr}(\hat{\mathbf{P}}\mathbf{S}\mathbf{S}^T)]$$

and

$$\xi = 2 \left[ \text{tr}(\hat{\mathbf{P}}\mathbf{S}\mathbf{S}^T) \right]^2 / \mathcal{I}_{\tau}.$$

$$\boxed{\frac{U_{\tau}}{\kappa} \sim \chi_{\xi}^2}$$



Sección 5

# Simulations

Simulated model

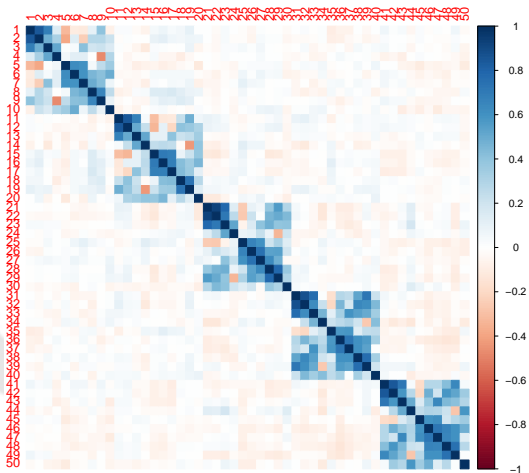
$$E_{ij} = 2 + 0,01Age_{ij} + 0,1I(Female_{ij}) + \gamma_i + \varepsilon_{ij};$$

$$\begin{aligned} & \text{logit} [P(Y_{ij} = 1 | Age_{ij}, Female_{ij}, E_{ij}, G1_{ij}, G2_{ij})] \\ = & 0,1 + 0,01Age_{ij} + 0,1I(Female_{ij}) + 0,1E_{ij} + 0,3G1_{ij} \\ & + 0,3G2_{ij} + \gamma_1(G1_{ij} \times E_{ij}) + \gamma_2(G2_{ij} \times E_{ij}) + \alpha_{ij} \end{aligned}$$

with

- $I(\cdot)$  is the indicator function;
- $\varepsilon_i = (\varepsilon_{i1}, \dots, \varepsilon_{i10})^T \sim N(\mathbf{0}, 4\mathbf{I}_{10})$ ;
- $\alpha_i = (\alpha_{i1}, \dots, \alpha_{i10})^T \sim N(\mathbf{0}, 2\Phi_i)$ ;
- $\gamma_i \sim N(0, 4)$
- $\Phi_i$  is the kinship matrix corresponding to the aforementioned family pedigree.

## Correlation of the 50 simulated SNPs in LD



## Type I error

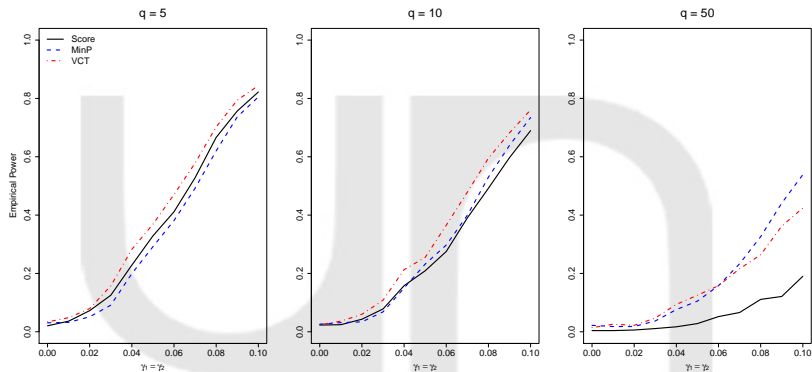
We first compared the empirical type I error of the different methods at 0.05  $\alpha$ -level. To evaluate type I error, we set  $\gamma_1 = \gamma_2 = 0$  and varied the number of SNPs  $q$  in the gene. The SNPs were either independent or in LD.

## Type I error

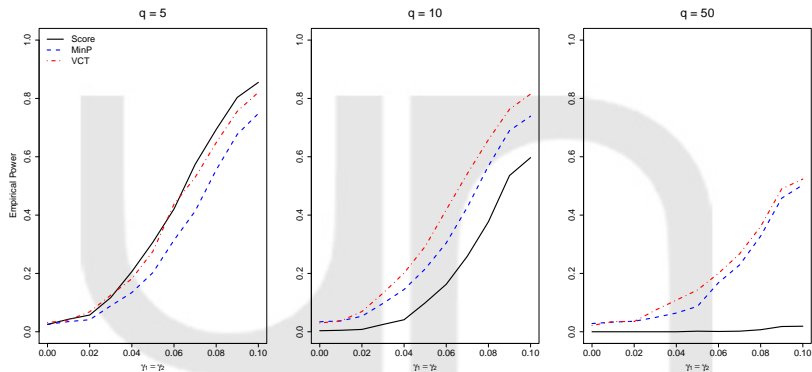
We first compared the empirical type I error of the different methods at 0.05  $\alpha$ -level. To evaluate type I error, we set  $\gamma_1 = \gamma_2 = 0$  and varied the number of SNPs  $q$  in the gene. The SNPs were either independent or in LD.

SNPs category	$q$	$\hat{\sigma}^2$	$\hat{\sigma}_\theta^2$	$\hat{\lambda} = 1/\hat{\sigma}_\theta^2$	Score	MinP	VCT
Independent	5	1.247	0.034	29.4	0.020	0.031	0.034
	10	1.240	0.017	58.8	0.023	0.026	0.024
	50	1.222	0.003	333	0.004	0.022	0.014
LD	5	1.243	0.021	47.6	0.025	0.026	0.031
	10	1.239	0.009	111	0.004	0.034	0.030
	50	1.222	0.002	500	0.000	0.028	0.022

## Empirical power for independent SNPs



## Empirical power for dependent SNPs



Sección 6

# Application



## Baependi Heart Study



## Variables

- Phenotype: Type II diabetes
- Environmental variable: Body Mass Index (BMI)
- Genotype: Were considered the following three genetic regions:
  - Peroxisome-Proliferator-Activated Receptors gamma (PPARG) with 16 variants genotyped;
  - Fat Mass and Obesity associated protein (FTO) with 149 variants genotyped;
  - Cyclin-dependent kinase 5 regulatory subunit associated protein 1-like 1 (CDKAL1) with 186 variants genotyped.
- Covariates: Age, sex, and the first two principal components of the entire genotype data of Baependi.

Summary of cases per subjects and families

Gene	Subjects			Families		
	Control	Cases	Total	Control	Cases	Total
PPARG	845	83	928	43	42	85
FTO	712	71	783	47	38	85
CDKAL1	661	69	730	47	38	85

## Summary of cases per subjects and families

Gene	Subjects			Families		
	Control	Cases	Total	Control	Cases	Total
PPARG	845	83	928	43	42	85
FTO	712	71	783	47	38	85
CDKAL1	661	69	730	47	38	85

R version 3.3.1 and a processor Intel(R) Core(TM) i5-6500 CPU @ 3.20GHz with a RAM memory 8.00 GB and operating system 64-bits.

Sample size, GLMM parameters, p-values and execution times

Gene	SNPs	Total Subjects	GLMM Parameters			Test	p-Value	$\alpha$ Level	Time (s)
			$\hat{\sigma}^2$	$\hat{\sigma}_\theta^2$	$\hat{\lambda} = 1/\hat{\sigma}_G^2$				
PPARG	16	928	0.4463	0.0029	344.8276	VCT	0.028	0.05	18.420
						MinP	0.019 *	0.005	100.261
						Score	0.595	0.05	9.025
FTO	149	783	0.3710	0.0033	303.0303	VCT	0.451	0.05	12.958
						MinP	0.031 *	0.0005	2675.907
						Score	0.992	0.05	6.197
CDKAL1	186	730	0.0918	0.0111	90.0901	VCT	0.635	0.05	9.907
						MinP	0.040 *	0.0005	1755.881
						Score	0.915	0.05	4.257

\* Compare MinP test  $p$ -value with the corresponding corrected  $\alpha$ , obtained by dividing 0.05 for the number of effective SNPs (which is equivalent to the number of principal components that reach 99.5% of the their total variation): 10 for *PPARG*, 93 for *FTO* and 92 for *CDKAL1*.

Sección 7

# References

References

- Breslow, N. and Clayton, D, 1993. Approximate inference in generalized linear mixed models. *J. Am. Stat. Assoc.*, **88**, 9-25.
- Coombes, B., Basu, S. and Mcgue, M. , 2017. A combination test for detection of gene-environment interaction in cohort studies. *Genet. Epidemiol.*, **41**, 396-412.
- Chen, H. and Conomos, M., 2016. GMMAT: Generalized Linear Mixed Model Association Tests; R Package Version 0.7. Available online:  
[https://content.sph.harvard.edu/xlin/dat/GMMAT\\_user\\_manual\\_v0.7.pdf](https://content.sph.harvard.edu/xlin/dat/GMMAT_user_manual_v0.7.pdf)  
(accessed on 16 January 2017).
- Lin, X., Lee, S., Chistiani, D. and Lin, X., 2013. Test for interactions between a genetic marker set and environment in generalized linear models. *Biostatistics*, **14**, 667-681.

References

- Lin, X., Lee, S., Wu, M., Wang, C., Chen, H., Li, Z. and Lin, X., 2016. Test for rare variants by environment interactions in sequencing association studies. *Biometrics*, **72**, 156-164.
- Lin, X., 1997. Variance component testing in generalised linear models with random effects. *Biometrika*, **84**, 309-326.
- Oliveira, C., Pereira, A., de Andrade, M., Soler, J. and Krieger, J., 2008. Heritability of cardiovascular risk factors in a brazilian population: Baependi heart study. *BMC Med. Genet.*, **9**, 32.
- Shen, X., Alam, M., Fikse, F. and Ronnegard, L., 2013. A novel generalized ridge regression method for quantitative genetics. *Genetics*, **193**, 1255-1268.
- Zhang, D. and Lin, X., 2003. Hypothesis testing in semiparametric additive mixed models. *Biostatistics* 2003, **4**, 7-74.



References

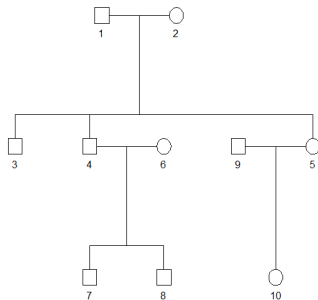
- Satterthwaite, F., 1941. Synthesis of variance. *Psychometrika*, **6**, 309-316.
- Fu, W.J., 2005. Nonlinear GCV and quasi-GCV for shrinkage models. *Journal of Statistical Planning and Inference*, **131**, 333-347.
- Mazo, M. A., Coombes, B. and de Andrade, M., 2017. An Efficient Test for Gene-Environment Interaction in Generalized Linear Mixed Models with Family Data. *International Journal Of Environmental Research And Public Health*, **14**, 1134-1146.

References

- Chen, H., Wang, C., Conomos, M., Stilp, A., Li, Z., Sofer, T., Szpiro, A., Chen, W., Brehm, J., Cedeon, J., Redline, S., Papanicolaou, G., Thornton, T., Laurie, C., Rice, K. and Lin, X., 2016. Control for population structure and relatedness for binary traits in genetic association studies via logistic mixed models. *The American journal of human genetics*, **98**, 653-666.
- Leal, S., Yan, K. and Muller-Myhsokb, B, 2005. SimPed: A Simulation Program to Generate Haplotype and Genotype 308 Data for Pedigree Structures. *Hum Hered.*, 119-122.

## Family data and kinship matrix

**PEDIGREE**



**KINSHIP MATRIX**

subject	1	2	3	4	5	6	7	8	9	10
1	1.00	0.00	0.50	0.50	0.50	0.0	0.250	0.250	0.0	0.250
2	0.00	1.00	0.50	0.50	0.50	0.0	0.250	0.250	0.0	0.250
3	0.50	0.50	1.00	0.50	0.50	0.0	0.250	0.250	0.0	0.250
4	0.50	0.50	0.50	1.00	0.50	0.0	0.500	0.500	0.0	0.250
5	0.50	0.50	0.50	0.50	1.00	0.0	0.250	0.250	0.0	0.500
6	0.00	0.00	0.00	0.00	0.00	1.0	0.500	0.500	0.0	0.000
7	0.25	0.25	0.25	0.50	0.25	0.5	1.000	0.500	0.0	0.125
8	0.25	0.25	0.25	0.50	0.25	0.5	0.500	1.000	0.0	0.125
9	0.00	0.00	0.00	0.00	0.00	0.0	0.000	0.000	1.0	0.500
10	0.25	0.25	0.25	0.25	0.50	0.0	0.125	0.125	0.5	1.000

Ridge regression estimation

$$g(\mu^{d_2}) = \tilde{\mathbf{X}}\beta + \mathbf{G}\theta + \mathbf{d}_2 \quad (\mathbf{d}_2 = \mathbf{K}\mathbf{b})$$



Ridge regression estimation

$$g(\mu^{d_2}) = \tilde{\mathbf{X}}\beta + \mathbf{G}\theta + \mathbf{d}_2 \quad (\mathbf{d}_2 = \mathbf{K}\mathbf{b})$$

$$\begin{aligned} ql(\beta, \theta, \phi_R, \sigma_R) &= -\frac{1}{2} \log \left| \frac{\sigma_R^2}{\phi_R} (\mathbf{K}\mathbf{K}^T) \mathbf{W}_R + \mathbf{I}_n \right| \\ &\quad + \sum_{i=1}^N \sum_{j=1}^{n_j} ql_{ij}(\beta, \theta, \phi_R; \tilde{\mathbf{d}}_2) - \frac{1}{2} \tilde{\mathbf{d}}_2^T (\sigma_R^2 \mathbf{K}\mathbf{K}^T)^{-1} \tilde{\mathbf{d}}_2, \end{aligned}$$

where  $\tilde{\mathbf{d}}_2$  is chosen to maximize the sum of the last two terms,  
 $\mathbf{W}_R = \text{diag} \left\{ \omega_{ij} / \left[ \nu(\mu_{ij}^{d_2}) g'(\mu_{ij}^{d_2})^2 \right] \right\}$  and

$$ql_{ij}(\beta, \theta, \phi_R; \mathbf{d}_2) = \int_{Y_{ij}}^{\mu_{ij}^{d_2}} \frac{\omega_{ij}(Y_{ij} - \mu)}{\phi_R \nu(\mu)} d\mu.$$

## Ridge regression estimation

Ridge regression estimator of  $\theta$ , is obtained by minimizing the function

$$[qI(\beta, \theta, \phi_R, \sigma_R) \times \phi_R] - \frac{1}{2} \lambda \theta^T \theta.$$

where  $\lambda$  is a penalizing factor.

$$\begin{bmatrix} \tilde{\mathbf{X}}^T \mathbf{W}_R \tilde{\mathbf{X}} & \tilde{\mathbf{X}}^T \mathbf{W}_R & \tilde{\mathbf{X}}^T \mathbf{W}_R \\ \mathbf{W}_R \tilde{\mathbf{X}} & \mathbf{W}_R + \lambda(\mathbf{G}\mathbf{G}^T)^{-1} & \mathbf{W}_R \\ \mathbf{W}_R \tilde{\mathbf{X}} & \mathbf{W}_R & \mathbf{W}_R + \frac{\phi_R}{\sigma_R^2}(\mathbf{K}\mathbf{K}^T)^{-1} \end{bmatrix} \begin{pmatrix} \beta \\ \mathbf{d}_1 \\ \mathbf{d}_2 \end{pmatrix} = \begin{pmatrix} \tilde{\mathbf{X}}^T \mathbf{W}_R \tilde{\mathbf{Y}} \\ \mathbf{W}_R \tilde{\mathbf{Y}} \\ \mathbf{W}_R \tilde{\mathbf{Y}} \end{pmatrix}$$

where  $\mathbf{d}_1 = \mathbf{G}\theta$ . [Return](#)